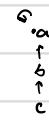
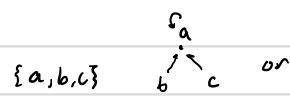
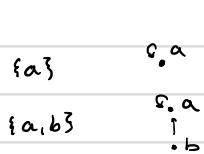


Can represent a set w/ a tree



$\pi(a) = a \therefore a$ is root
 $\pi(b) = a$
 $\pi(c) = b$
 $\text{rank}(f) = 0$
 $\text{rank}(b) = 1$
 $\text{rank}(a) = 2$

Defn: $\pi(x)$ = parent node (think "pointer")

root node = the x s.t. $\pi(x) = x$

use the root node to ID sets

$\text{rank}(x)$ = # edges in longest simple path from x to a leaf node

$\therefore \text{def make_set}(x)$

$\pi(x) = x$

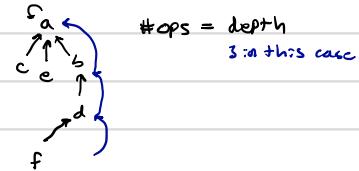
Set parent to self

$O(1)$

$\text{rank}(x) = 0$

~~def find(x)~~
while $x \neq \pi(x)$
| $x = \pi(x)$
return x

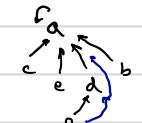
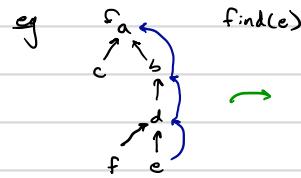
look for root node



~~def find(x)~~
while $x \neq \pi(x)$
| $\pi(x) = \text{find}(\pi(x))$
return $\pi(x)$

Path compression

Two pass method:
1. up to find root
2. make each node on
found path point to
root node



$\text{find}(e)$

Subsequent finds() on this tree $O(1)$

requires amortized analysis
i.e. over a sequence of operations

But what about union?

thus

$\text{def union}(x, y)$ ↗ aka, Union by rank

$r_x = \text{find}(x)$

$r_y = \text{find}(y)$

If $r_x = r_y$

| return

If $\text{rank}(r_x) > \text{rank}(r_y)$

| $\pi(r_y) = r_x$

Else

| $\pi(r_x) = r_y$

| If $\text{rank}(r_x) = \text{rank}(r_y)$

| | $\text{rank}(r_y) += 1$

$$a \cup b = \begin{array}{c} a \\ | \\ b \end{array}$$

Note: attach smaller ranked tree
to larger one so $\text{find}(.)$ won't take
too long

Kruskal running time w/ tree data structure

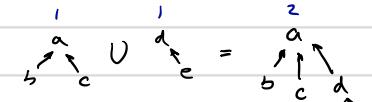
```

def Kruskal-MST(G, w)
    A = ∅
    for all v ∈ G.V
        | make-set(v) ] O(|V|)

    G.E = G.E Sorted in non-decreasing order by w ↪ O(E log E) = O(E log V) ← since E ≤ V2
    for all (u,v) ∈ G.E
        | if find-set(u) ≠ find-set(v)
        |   | A = A ∪ {(u,v)} ← 2nd loop
        |   | union(u,v)
    return A

```

- Observe:
1. $\forall x, \text{rank}(x) < \text{rank}(\pi(x))$
 2. Root node w/ rank K is formed from the union of 2 rank $K-1$ trees
 3. Any root node of rank K has at least 2^K nodes in the tree



Proof: Induction $K=0 \quad \exists a \quad 2^0=1 \checkmark$

Assume true @ $K-1$. By prop 1, need to union 2 trees of rank $K-1$, each has 2^{K-1} nodes
 \Rightarrow tree w/ root of rank K has $2^{K-1} + 2^{K-1} = 2^K$ nodes \square

4. w/ $|V|$ nodes, max rank is $\log_2 |V|$

$\therefore \text{find}(\cdot) = O(\log |V|)$

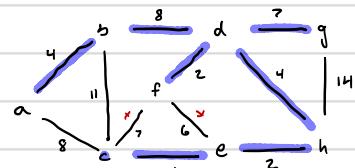
$\text{union}(\cdot) = O(\log |V|) \leftarrow 2 \text{ finds + constant work} \quad \therefore \text{2nd loop } O(|E| \log |V|)$

So total running time $O(|E| \log |V|)$

Prim's algorithm

Similar, but iteratively grows tree

Form cut w/ nodes used so far, pick lightest edge crossing that cut



See text for more details

Huffman encoding

Used in MP3 compression

Data: string S over an alphabet Γ $\xrightarrow{\text{very long}} \text{usually } |\Gamma| > 2$

Goal: Find a binary encoding of Γ resulting in a minimal encoding length S_e

$$\text{eg } \Gamma = \{a, b, c\}$$

Statistics of S : # char

a b c
45 16 2

$$\text{fixed-length encoding: } a \rightarrow 00 \quad |S_e| = 2 \cdot 45 + 2 \cdot 16 + 2 \cdot 2 = 126$$

$b \rightarrow 01$
 $c \rightarrow 10$

Wasted space... what about variable length encoding?

$$a \rightarrow 0 \quad T_e = 1 \cdot 45 + 2 \cdot 16 + 2 \cdot 2 = 81$$

$b \rightarrow 10$
 $c \rightarrow 11$

Careful about ambiguity! If code words were $a \rightarrow 0$, $b \rightarrow 1$, $c \rightarrow 01$ then 01 is ambiguous $01 = c?$
 $01 = ab?$

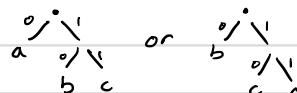
Need our encoding to be prefix-free: no code word is prefix of another code word

Use a full binary tree: each node has 0 or 2 children

left edge: label 0

right edge: label 1

decorate leaves w/ letters from Γ



$b \rightarrow 0$
 $c \rightarrow 10$
 $a \rightarrow 11$

encoding reads off edge labels on path from root to leaf

guaranteed to be prefix-free

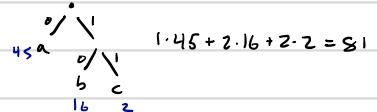
$\nwarrow \text{depth}(b) = 1$

$\text{depth}(c) = \text{depth}(a) = 2$

$\downarrow \text{think "frequency"}$

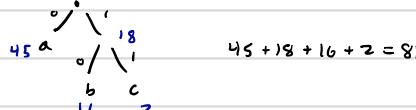
So given encoding tree T , string S w/ letter counts f_r , Γ , size of encoded string is

$$|S_e| = \sum_{v \in T} f_r \cdot \text{depth}(v) \quad \leftarrow \text{want to minimize this wrt encoding tree } T$$



Clearer re-write: label internal nodes w/ sum of descendant counts

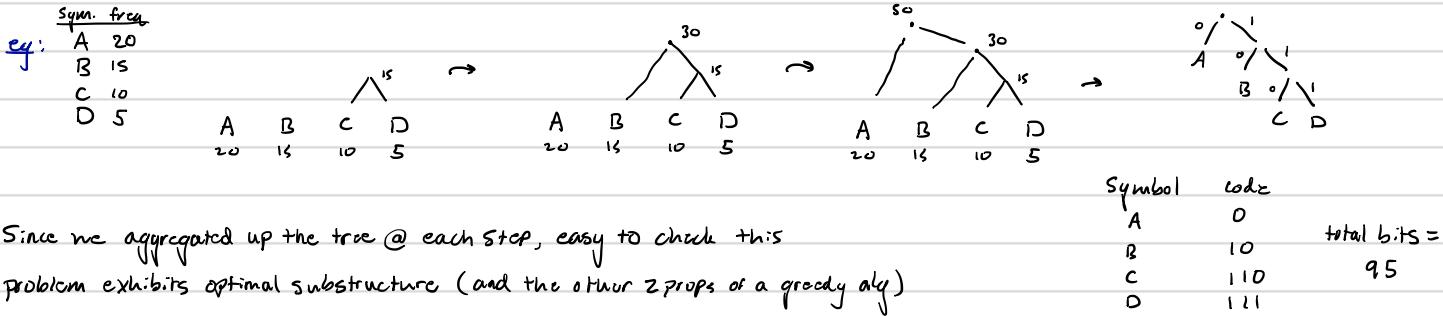
$\forall v \in T$, let $\text{cost}(v) = \text{sum of leaf node counts descending from } v$



$$\text{then } |S_e| = \sum_{v \in T} \text{cost}(v)$$

This suggests a greedy approach: choose encoding of length inversely proportional to symbol frequency (i.e. most frequent symbols have smaller depth)

Idea behind algorithm: Continually merge least frequent nodes until you have a binary tree



Since we aggregated up the tree @ each step, easy to check this problem exhibits optimal substructure (and the other 2 props of a greedy alg)

def Huffman(f):
 array f[1...n] of frequencies

Let T be an empty tree $O(1)$

Let H be a priority queue, ordered by f $O(1)$

for i=1 to n $n \cdot O(H \text{ insert time})$

 insert(H, i)

for K=n+1 to 2n-1

 i = extract min(H)

 j = extract min(H)

$(n-1)(2 \cdot O(H \text{ extract min}) + O(1) + O(1) + O(H \text{ insert time}))$

 put node K in T with children i and j

 f[K] = f[i] + f[j]

 insert(H, K)

 insert, min, extract, decrease

So complexity depends on implementation of the priority queue

Hcap-based: insert $O(1)$, extract min $O(\log n) \Rightarrow$ Huffman is $O(n \log n)$

van Emde Boas tree based: insert $O(\log \log n)$, extract min $O(\log \log n) \Rightarrow$ Huffman is $O(n \log \log n)$

Side note re: entropy

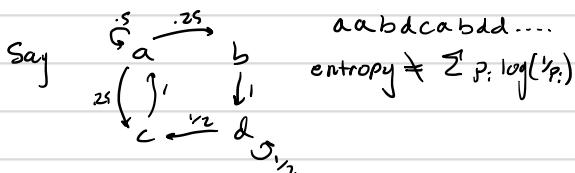
n symbols w/ freq $P_1, \dots, P_n \quad \sum_i P_i = 1$

pull m total values according to the dist (P_1, \dots, P_n)

are number of bits to encode (using Huffman): $\sum_{i=1}^n m \cdot P_i \log(\frac{1}{P_i}) \rightarrow$ single draw $H(P_1, \dots, P_n) := \sum_{i=1}^n P_i \log(\frac{1}{P_i})$

Only applicable when iid! (one @ a time, indep, all pulled according to same dist).

Relatedly, Huffman only optimal in this situation \rightarrow w/ restriction of 1 symbol \rightarrow 1 code word



metric / Shannon's entropy

Matroids, processor scheduling

Set cover

DP

LP

NP

Matroids

Theory generalizing many greedy algorithms

Def'n a matroid $M = (S, \mathcal{I})$ is an ordered pair s.t.

1. S is a finite set

2. $\mathcal{I} \subseteq \mathcal{P}(S)$ set of subsets (call them indep sets) s.t.

$B \in \mathcal{I}$ and $A \subseteq B$, then $A \in \mathcal{I}$

all subsets of each set included

3. exchange property: $A \neq B \in \mathcal{I}$ s.t. $|A| < |B|$, then $\exists x \in B - A$ s.t. $A \cup \{x\} \in \mathcal{I}$

e.g. Graphic Matroid

given undirected graph $G = (V, E)$, $M_G = (S_G, \mathcal{I}_G)$ defined via

$S_G := E$ (set of edges)

$A (\subseteq E)$ is a $\tilde{\Delta}_G$ if A is acyclic (i.e. if $G_A = (V, A)$ is a forest)

$$G = \boxed{1 \quad 2 \quad 3 \quad 4 \quad 5}$$

$$S = \{(1,2), (1,3), (2,3), \dots, (4,5)\}$$

$\mathcal{I}_G = \text{set of all forests of } G$

$$= \{ \{(1,2)\}, \{(1,2), (1,3)\}, \{(1,2), (1,3), (3,4)\}, \dots, \}$$

\uparrow
 $\{(1,3)\}$ is an extension
of $\{(1,2)\}$

Can verify this is a matroid

Def'n For $A \in \mathcal{I}$, say $x \in S$ is an extension of A if $A \cup \{x\} \in \mathcal{I}$

Def'n $A \in \mathcal{I}$ is maximal if it has no extension

Thm All maximal $A \in \mathcal{I}$ have the same size

Proof isoac, say $\exists B \in \mathcal{I}$ s.t. $|B| > |A|$. Then exchange property says $\exists x \in B - A$ s.t. $A \cup \{x\} \in \mathcal{I}$ \Rightarrow contradicts maximality of A . \square

Def'n A weighted matroid is a matroid $M = (S, \mathcal{I})$ along w/ Strictly positive $w: S \rightarrow \mathbb{R}_{>0}$

extend to sets via: $A \subseteq S \Rightarrow w(A) = \sum_{x \in A} w(x)$

e.g. In graph matroid, w gives edge weights

Some const larger than 2

Observe: minimum spanning tree can be cast as a "maximum indep set" problem by setting $w(\text{edge}) = C - \text{orig.-edge-weight}(\text{edge})$

So if we could find an alg. that solves the "find a maximum indep set" for matroids, we would have another MST algorithm!

def Greedy($M = (S, \mathcal{I}), w$)

$A = \emptyset$

let $n = |S|$

$S = \text{Sort } S \text{ into decreasing order by } w$

$O(n \log n)$

for each $x \in S$

Say checking $A \cup \{x\} \in \mathcal{I}$ is $O(f(w))$

if $A \cup \{x\} \in \mathcal{I}$
 $A = A \cup \{x\}$

Oracle

return A

} total time:

$O(n \log n + n f(w))$

Can show this exhibits all 3 prop (iterative, safe, optimal substructure) and so returns an optimal/maximal A.

Task scheduling as a matroid problem

e.g. CPU scheduling programs to run

Setup: $S = \{a_1, \dots, a_n\}$ n unit-time tasks

d_1, \dots, d_n deadlines by which tasks must be completed, $1 \leq d_i \leq n$

w_1, \dots, w_n non-negative penalties (weights) we incur if a_i is not finished by d_i deadline

Schedule: perm. of S specifying when to do tasks late early still early still late

Can transform any schedule into early first: $\{ \dots, a_i, \dots, a_j, \dots \} \rightarrow \{ \dots, a_j, \dots, a_i, \dots \}$

Canonical form: early before late, early ordered w/ monotonically increasing deadlines (snap early if not)

Defin a set A of tasks is indep. if \exists a way to schedule them so that no task is late

Let \mathcal{I} = set of all indep. tasks

Clearly \mathcal{I} is hereditary, it possesses the exchange property: if $|B| > |A|$, $B, A \in \mathcal{I}$, let K be the time when last task in A finishes, let x be the first task in B that finishes after K , then $A \cup \{x\} \in \mathcal{I}$.

$\therefore (S, \mathcal{I})$ is a matroid and so \exists greedy alg. to solve it $\rightarrow O(n \log n + n^2) = O(n^2)$

checking $A \cup \{x\} \in \mathcal{I}$ takes $O(n)$

Last couple of algs

Horn Formulas

Boolean Formulas: $x = T \text{ or } F$, $x \wedge \bar{y} \neq z \vee a$
 and implication

Special Kinds: Horn Formulas: a list of Boolean formulas of the form:

1. implications: $(x_1 \wedge \dots \wedge x_n) \Rightarrow y$ all pos. literals, "ands" on LHS, Single literal on the right
 2. negative clauses: $(\bar{x}_1 \vee \dots \vee \bar{x}_n)$ all neg. literals, "ors" only

Goal: Given a set of Horn formulas, find values of all literals that makes it evaluate to T (or report if you can't)

$$\text{eg } (x \wedge y) \Rightarrow z, (\bar{x} \vee \bar{w}), \quad \text{Say } \begin{array}{l} w = F \\ x = F \\ y = F \\ z = F \end{array} \quad \text{counter eg } \begin{array}{l} x \neq y, \quad \bar{x}, \quad \neq y \\ F \neq T \quad x = F \quad y = T \end{array}$$

Idea: Start w/ all F, only set T if you need to, i.e. when $a_n \Rightarrow$ says you need to

Recall: $p \Rightarrow q \equiv \neg p \vee q$

P	q	$P \neq q$
T	T	T
T	F	F
F	T	T
F	F	T

→ means we need to set $q = T$ to get T

Correctness: Thus Any literal set to T by Greedy Horn
must be set to true in any sat. assignment

def Greedy Horn (set of Horn formulas)

Set all literals to F

while $a_n \Rightarrow$ is not satisfied

Set RHS to T
all negative clauses 0
return assignment

else

return "cannot be satisfied

$O(n^2)$ but can be reduced to $O(n)$

$$\text{eg } (x \wedge y) \Rightarrow z, (\bar{x} \vee \bar{w}), \quad x \quad y \quad z \quad (x \wedge y) \Rightarrow z \quad (\bar{x} \vee \bar{w})$$

$\top \quad \top \quad \top \quad \top \quad \top$

$$\begin{array}{ll} \text{eg} & \Rightarrow x, (\bar{x} \vee \bar{y}), x \Rightarrow y \\ & \begin{array}{ccccc} x & y & \Rightarrow x & (\bar{x} \vee \bar{y}) & x \Rightarrow y \\ \top & \top & \top & \top & \top \\ \top & \bot & \top & \top & \top \\ \bot & \top & \top & \top & \top \\ \bot & \bot & \top & \top & \top \end{array} \end{array}$$

\therefore can't satisfy

Set Cover Resource allocation, code testing, shift planning, etc

input: 1. Set B

2. Subsets $S_1, \dots, S_n \subseteq B$

output: a selection of subsets $I = \{S_{i_1}, \dots, S_{i_m}\}$ s.t. $\bigcup_{j=1}^m S_{i_j} = B$

Want to minimize $|I|$, the number of subsets selected

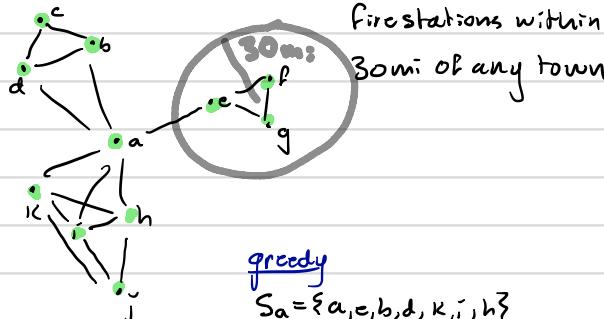
$$\text{eg } B = \{a, b, c, \dots, k\}$$

$$S_a = \{a, b, d, e, h, i, k\}$$

$$S_b = \{b, a, c, d\}$$

$$S_k = \{k, a, h, i, j\}$$

i.e. the union of the Selected subsets is B



greedy heuristic: Choose the next subset based on the most # of uncovered cities

not optimal!

$$S_b, S_c, S_i$$

But! Not off by too much

Thm Say $|B|=n$, optimal soln uses K sets, then greedy heur. uses $\leq K \ln(n)$ sets

$$\text{eg } K=3 \\ n=11$$

Proof: Let $n_t = \#$ elements not covered by greedy @ iteration t

these remaining elems are covered by an optimal soln, so some optimal

Set has size $\geq \frac{n_t}{K}$ elements of B in it. \therefore

called the approximation factor

$$3 \ln(11) = 7.19$$

$$n_{t+1} \leq n_t - \frac{n_t}{K} = n_t \left(1 - \frac{1}{K}\right)$$

greedy picks set
@ least this size

$$\text{repeat: } n_{t+1} \leq n_t \left(1 - \frac{1}{K}\right) \leq n_{t-1} \left(1 - \frac{1}{K}\right)^2 \leq \dots \leq n_0 \left(1 - \frac{1}{K}\right)^t = n \left(1 - \frac{1}{K}\right)^t \leq n \left(e^{-\frac{1}{K}}\right)^t = n e^{-t/K}$$

$$n_t \leq n_{t-1} \left(1 - \frac{1}{K}\right)$$

general fact via calculus:

$$1-x \leq e^{-x}$$

$$\therefore 1 - \frac{1}{K} \leq e^{-\frac{1}{K}}$$

Now if $n_t < 1$, the program terminates, so want to solve $n_t < 1$ for t . Since $n_t < n e^{-t/K}$ if $n e^{-t/K} < 1$
 $\therefore n e^{-t/K} = 1 \Leftrightarrow e^{-t/K} = \frac{1}{n} \Leftrightarrow -t/K = \ln(1/n) \Leftrightarrow t = -K \ln(1/n) = K \ln(n)$.
 \therefore will terminate within $t = K \ln(n)$ iterations. \square

Dynamic Programming

In greedy methods, making the greedy choice meant only one subproblem would remain.

If this is not the case (not 1 problem optimal substructure) will have multiple subproblems to analyze.

→ dynamic programming is designed for cases like this

hence: more generally applicable, but less efficient

Applies to optimization problems

General Steps

Break the problem into smaller subproblems

tackle subproblems smaller first (bottom up typically)

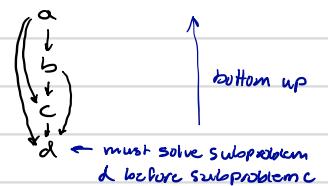
use information from smaller problems to solve bigger ones

This is conceptualized with a subproblem graph: a DAG where vertices represent subproblems, edges $x \rightarrow y$ if determining an optimal solution for x directly involves considering an optimal solution for y .

Note: typically, time to solve a subproblem is proportional to # of outgoing edges

Subproblems = # vertices

So \propto running time linear in # vertices + # edges



e.g. Longest incr. subsequence

Set of #'s a_1, \dots, a_n , find a subsequence $a_{i_1}, a_{i_2}, \dots, a_{i_k}$ w/ $i_1 < i_2 < \dots < i_k$ w/ largest length K

$S \ 2 \ 8 \ 6 \ 3 \ 6 \ 9 \ 7$

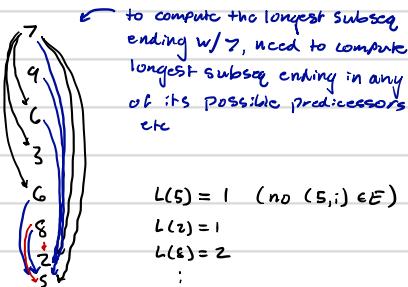
Subproblem DAG!

Solve from the bottom up!

For v in reverse topological sort of $(V, E) \rightarrow$

$$L(v) = 1 + \max \{ L(i) : (v, i) \in E \}$$

(Return $\max_v L(v)$)



Designing a DP

look for

1. optimal substructure

2. overlapping subproblems

In general, for $P(i, j)$, need only pick which of 3 cases is least cost:

$$P(i, j) = \min (1 + P(i-1, j), 1 + P(i, j-1), d: P(i, j) + P(i-1, j-1))$$

$$\curvearrowleft = \begin{cases} 0 & x_i = y_j \\ 1 & x_i \neq y_j \end{cases}$$

∴

def EditDist(x, y):

for $i=0, \dots, m$

$$P(i, 0) = i$$

for $j=0, \dots, n$

$$P(0, j) = j$$

for $i=0, \dots, m$ ↓ order doesn't matter

for $j=0, \dots, n$

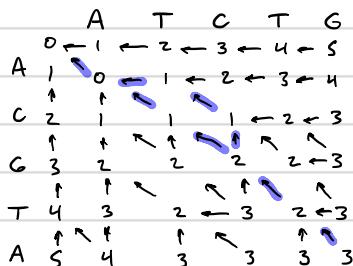
$$P(i, j) = \min (1 + P(i-1, j), 1 + P(i, j-1), d: P(i, j) + P(i-1, j-1))$$

return $P(m, n)$

$O(mn)$

got actual alignment by keeping track
of where min came from & backtracking
from $P(m, n)$

eg



$$\therefore P(m, n) = d_E(\text{ACGTA}, \text{ATCTG}) = 3$$

↑ = match

↑ = blank in top of alignment

↔ = blank in bottom "

A C G T A
A T C T G

or
ATC-TA
A-CGTG

cost 3

cost 3

Weighting ↑ or ↓ or ↔ by different amounts leads to alternative alg's: Smith-Waterman, Needleman-Wunsch

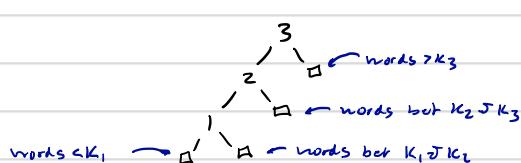
example 2: instead of chain matrix mult (see book), problem in similar spirit: optimal binary search

Translate French book to English, organize dict. for efficiency: word freq is Zipfian ↓ takes word freq into account

Some words may have no translation

Say only 3 words w/ keys $k_1 < k_2 < k_3$ w/ freq/prob p, q, r

Say lexicographic order



$$\text{cost } 3p + 2q + r$$



etc

$$\text{cost } p + 3q + r$$

Idea is most frequent words

Should be highest in tree: cost
= total tree traversal length over
whole corpus

$\exists \binom{n}{n+1} \approx 4^n / (7 \cdot n^{3/2})$ binary trees on n nodes → can't brute force check for lowest cost tree

Also, absent word freqs will play a role (the □ in the above).